

第七章 分类 读书笔记

7.1 分类算法概要

7.1.1 分类定义及原理

分类是数据挖掘中的一个重要任务，它通过分析数据特征来预测样本所属的类别。对于数据集 $D=\{x_1, x_2, \dots, x_n\}$ ，每个样本 $x=(x_1, x_2, \dots, x_d, c)$ 具有 d 个特征和一个类别标签 c 。分类的目标是建立一个模型，能够根据样本的特征准确地预测其类别标签。在实际应用中，分类任务可以分为二分类、多分类和多标签分类。二分类任务预测样本属于两个类别中的一个，多分类任务预测样本属于多个类别中的一个，而多标签分类任务则为每个样本预测多个类别标签。

分类算法的核心在于如何确定决策边界。决策边界是在特征空间中将不同类别的数据样本分隔开的线、曲线或超平面。通过计算决策边界，模型能够判断新样本的类别。作为监督学习的一种，分类算法需要利用训练集中的样本特征及其类别标签来训练模型，以便对新的未标记数据进行准确预测。掌握分类算法对于理解数据背后的模式和规律至关重要，它能够帮助我们从大量数据中提取有用的信息，为决策提供依据。

7.1.2 常用分类算法

分类算法可以根据其理论基础进行分类。基于统计学原理的分类算法包括朴素贝叶斯和Logistic回归；基于机器学习的分类算法包括K-近邻、决策树和支持向量机；基于深度学习的分类算法则包括神经网络算法。

- **朴素贝叶斯**：基于贝叶斯定理和特征条件独立性假设，通过计算给定样本属于各个类别的概率来实现分类。它简单易懂，适用于处理文本数据等特征之间相对独立的情况。但在实际应用中，特征之间往往存在相关性，这可能会影响朴素贝叶斯算法的性能。
- **Logistic回归**：根据训练集建立回归模型来描述分类决策边界，并以此进行分类，通常用于二分类问题。它能够输出概率值，方便我们进行概率解释和决策。然而，在处理非线性关系时，Logistic回归可能不如其他算法表现突出。
- **K-近邻**：基于距离度量找出待预测样本与训练集中距离最近的 K 个训练样本，然后参考这 K 个“邻居”样本的类别来推断给定样本的类别。K-近邻算法简单直观，适用于对数据分布没有太多先验知识的情况。但在大规模数据集上，计算距离可能会非常耗时。
- **决策树**：通过训练集样本计算不同特征的分类能力来构建模型，从而对未知的数据进行分类。决策树易于理解和解释，能够处理数值型和分类型特征。不过，决策树容易过拟合，需要进行剪枝等操作来提高泛化能力。
- **支持向量机**：把分类问题转化为寻找分类超平面的问题，并通过最大化分类边界点到分类平面的距离来实现分类。支持向量机在处理高维数据和非线性问题时表现出色，但参数选择和核函数的选取对模型性能有很大影响。
- **神经网络**：由输入层、隐藏层和输出层组成的多层模型，通过学习将输入映射到输出的函数参数实现分类。神经网络具有强大的非线性映射能力，能够处理复杂的分类任务。然而，神经网络的训练过程复杂，需要大量的计算资源和时间。

每种算法都有其优势和局限性。选择合适的分类算法需要根据具体问题的特点和数据集的性质来决定。例如，在处理文本分类问题时，朴素贝叶斯算法可能是一个不错的选择；而在处理图像识别等复杂问题时，神经网络可能更具优势。此外，算法的性能也受到数据质量、特征选择和模型参数等因素的影响。因此，在实际应用中，我们通常需要对算法进行调优和优化，以获得更好的分类效果。

7.1.3 分类算法评价指标

为了评估和比较不同分类算法的性能，我们使用一系列评价指标。这些指标能够帮助我们了解模型在不同方面的表现，从而选择最适合的算法。

- **混淆矩阵**：用于比较模型的预测结果与真实情况。混淆矩阵是一个二维矩阵，行代表真实标签，列代表预测标签。对于二分类问题，混淆矩阵可以表示为：
 - 真正 (True Positive, TP)：预测为正类，实际为正类；
 - 真负 (True Negative, TN)：预测为负类，实际为负类；
 - 假正 (False Positive, FP)：预测为正类，实际为负类；
 - 假负 (False Negative, FN)：预测为负类，实际为正类。
- **准确率 (Accuracy)**：表示被正确分类的样本数占总样本数的比例。准确率能够反映模型对所有样本进行正确分类的概率，但在样本不均衡的情况下，准确率可能无法准确衡量模型的性能。例如，在一个正负样本比例悬殊的数据集中，即使模型总是预测为多数类，准确率也可能很高，但这并不代表模型具有良好的分类能力。
- **精准率 (Precision)**：又叫查准率，是指预测为正样本的数据中真正为正样本的比例。精准率适用于希望使模型的预测结果尽可能不出错的场景。例如，在医疗诊断中，我们更关注精准率，以避免误诊。
- **召回率 (Recall)**：又叫查全率、灵敏度，是指实际为正样本的数据中被预测为正样本的比例。召回率衡量了模型将所有正样本成功识别出的比例，适用于希望尽可能全地预测出正样本的场景。例如，在垃圾邮件过滤中，我们更关注召回率，以避免漏掉垃圾邮件。
- **特异性 (Specificity)**：又叫真负率，是模型在所有真实负样本中成功预测出的比例。特异性衡量了模型将所有负样本成功识别出的比例。在某些应用中，如疾病筛查，特异性也是一个重要的指标。
- **F1得分**：将精准率和召回率合并为一个指标，以便模型之间性能的比较。F1得分同时考虑了模型在正样本和负样本上的表现，适用于数据集类别分布不平衡的分类任务。
- **AUC (Area Under Curve)**：ROC曲线下的面积，用于评估模型在不同分类阈值下的性能。AUC的取值范围在0到1之间，当AUC=0.5时，模型的预测效果接近随机猜测；当AUC接近1时，模型的预测效果较好。AUC是一个重要的指标，尤其在处理不平衡数据集时，它能够提供更全面的性能评估。
- **马修斯相关系数 (Matthews Correlation Coefficient, MCC)**：考虑了真正例、真负例、假正例和假负例的二分类模型性能评估指标。MCC的取值介于-1和1之间，其中1表示完美的分类器，0表示随机分类器，-1表示完全相反的分类器。MCC适用于不平衡数据集的模型评价，能够提供比其他指标更平衡的性能评估。

单一指标往往无法全面反映模型的性能。例如，准确率可能会受到样本分布的影响，而精准率和召回率则需要根据具体应用场景来权衡。因此，在评估分类模型时，我们通常会综合考虑多个指标，以获得更全面的性能评估。此外，选择合适的评价指标也需要根据问题的特点和需求来决定。例如，在某些应用中，我们可能更关注模型的精准率，而在其他应用中，召回率可能更为重要。通过合理使用这些评价指标，我们可以更好地理解 and 比较不同分类算法的性能，为实际应用提供指导。

7.2 K-近邻

K-近邻 (K-Nearest Neighbor, KNN) 是一种经典的监督学习算法，其核心思想可概括为“物以类聚，人以群分”。它没有显式的学习训练过程，适用于对数据的分布只有很少甚至没有任何先验知识的情况。K-近邻算法的实现非常简单直观，其基本原理是基于特征空间中的距离度量，通过找到待分类样本的K个最近邻居（即特征空间中距离最小的样本），来确定其所属的类别。

在学习K-近邻算法时，我感受到了它独特的魅力。它不需要复杂的模型构建过程，只需在预测时根据最近的邻居样本进行决策，这种“懒惰学习”的方式使得算法的实现变得简单易懂。同时，K-近邻算法不依赖于数据的分布假设，能够适应各种类型的数据集，这使得它在处理一些复杂或未知分布的数据时具有一定的优势。

7.2.1 K-近邻基本原理

K-近邻算法的计算步骤如下：

1. 基于距离计算公式（如欧氏距离、曼哈顿距离等）计算待预测样本与每个训练样本之间的距离。
2. 对距离进行升序排序。
3. 筛选出与该测试样本距离最近的K个训练样本。
4. 将K个样本所属类别的频数与决定权重相乘，将乘积最高的类别作为测试样本的类别。

以图7-2为例，根据K-近邻的工作原理，当K=3时，与待分类样本最相近的是2个A类样本和1个B类样本，根据多数投票规则预测待分类样本的类别为A。当K=7时，与待分类样本距离最近的样本为3个A类样本和4个B类，则预测待分类样本的类别为B。

在理解K-近邻的基本原理后，我意识到算法的性能在很大程度上取决于K值的选择以及距离度量方法的合理性。K值的大小直接影响了决策边界，进而影响分类结果的准确性。同时，距离度量方法的选择也会影响样本之间的相似度计算，从而对分类结果产生影响。因此，在实际应用中，我们需要根据具体问题和数据集的特点来合理选择K值和距离度量方法

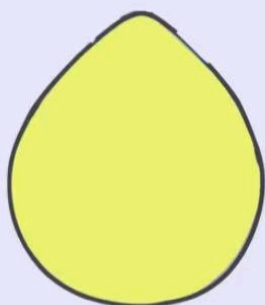
k-近邻算法

KNN - K nearest neighbour

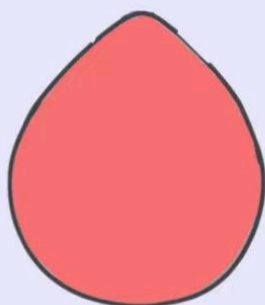
KNN (K-Nearest Neighbors) 算法，又称 **K 近邻算法**，常用的**监督学习**算法，主要用于**分类**和**回归**任务。它通过计算新样本与训练集中所有样本的距离，找出距离最近的 **K 个邻居**，然后根据这 **K 个邻居** 的类别或数值，通过**投票**或**平均**来预测新样本的类别或数值



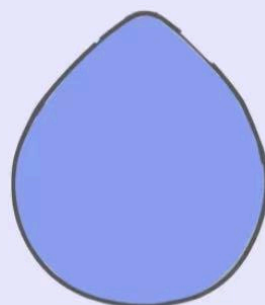
你们三个一会儿会降落到不同地方，刚到地球的任务，就是自己去寻找合适的身份，融入人类中



怎么找？



按剧本来讲，这篇是KNN，我们肯定要用KNN的方法找



收到

k-近邻算法

KNN - K nearest neighbour

工作原理：

1. 计算距离： 给定一个新样本，计算它与训练集中所有样本的距离。常用的距离度量方法包括欧氏距离、曼哈顿距离等。



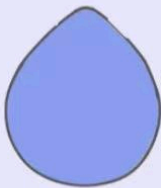
世界上醉遥远的距离不是生与
si，而是我就站在你面前…



闭嘴，飞太久已经很晕了



却不知道用什么度量方法测距离

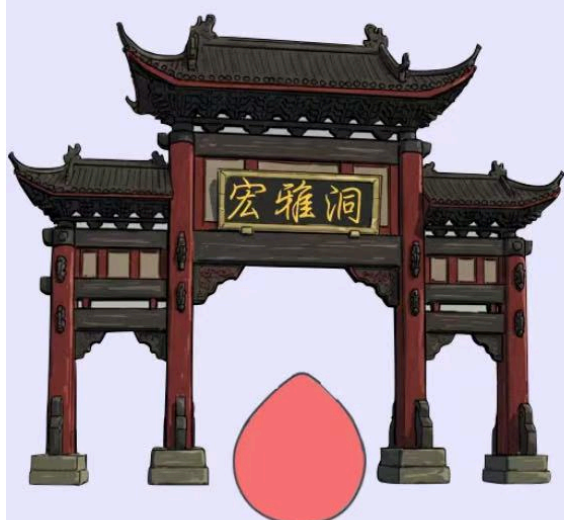


人类平常算距离用的就是欧氏距离，它表示的是
多维空间中两点之间的距离
一会儿到了我们要测量与所有人的距离
记得是所有，你们两个别偷懒

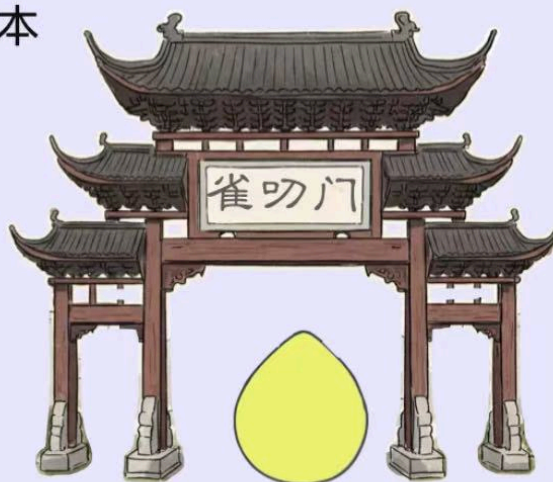
k-近邻算法

KNN - K nearest neighbour

2.寻找邻居： 选择距离最近的 K 个样本作为新样本的邻居。如k=5, 则选择最近5个样本



新邻居好像还不错



最近的k个都在这里了



k-近邻算法

KNN - K nearest neighbour

3. 预测类别/数值:

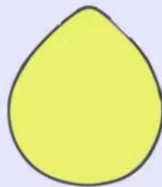
- **分类任务:** voting, K 个邻居中, 哪个类别出现次数最多, 就将新样本归为该类别。
- **回归任务:** averaging, 将 K 个邻居的数值取平均, 作为新样本的预测值。



找到自己的身份了吗?



是, 邻居们大部分是什么身份, 我们就是什么身份
绝对不会暴露



彪子



瓜娃子



赤佬



你们确定这样标记自己不会暴露吗...

7.2.2 K-近邻主要参数

7.2.2.1 K值

K值的选取对K-近邻算法性能至关重要。如果K值选取过小, 相当于用较少的训练样本进行预测, 只有与待分类样本非常相似的训练样本才会对预测结果起作用, 会导致待预测样本对近邻的训练样本十分敏感。如果训练集中存在噪声, 则极易导致预测出错。相反, 如果K值选取过大, 相当于用参考范围过大邻域内的样本作为决策依据, 会导致模型欠拟合。

- 从数据集大小的角度出发, 如果数据集较小, 选择较小的K值可能更合适, 通常K的取值一般不超过20, 或低于训练样本数N的平方根, 即 $K < \sqrt{N}$ 。

- 从数据特征的角度出发，如果数据集中样本的特征具有较大的方差或噪声，选择较大的K值更合适，因为较大的K值可以提供更平滑的决策边界。相反，如果数据集中的样本具有较小的方差，则较小的K值更合适。
- 从类别平衡性的角度出发，如果类别之间的样本数差异较大，较大的K值更合适。否则，选取较小的K值可能会导致预测结果被样本数较多的类别主导。

这是一个需要权衡的过程。K值过小可能导致模型过拟合，而K值过大则可能导致模型欠拟合。因此，在实际应用中，我们通常需要通过交叉验证等方法来确定最优的K值，以获得更好的分类性能。

7.2.2.2 距离度量方法

距离度量是指计算待预测样本与训练样本在特征空间上的距离。常用的距离度量方法包括欧氏距离、曼哈顿距离、闵可夫斯基距离、余弦距离等。为了避免特征量纲不一致对距离度量带来的误差，在计算距离之前通常会对特征做规范化处理。

不同的距离度量方法适用于不同类型的数据。例如，欧氏距离适用于连续型特征，而曼哈顿距离则适用于离散型特征。选择合适的距离度量方法需要根据数据的特征类型和分布情况来决定。此外，在处理具有不同量纲的特征时，特征规范化是必不可少的步骤，它能够确保各个特征在距离计算中具有相同的权重，从而提高分类结果的准确性。

7.2.2.3 权重函数

权重函数用于给K个最近邻居的投票进行加权，以确定待预测样本的类别。常用权重函数包括统一函数和距离权重函数。

- 在统一函数中，所有K个最近邻样本对最终决策的影响是相同的。对于待预测样本 x_g ，统计各类别在K个邻居样本中出现的频数，将频数最大的类别作为测试样本的类别。
- 在距离加权函数中，不同最近邻样本对最终决策的影响是根据其与待预测样本 x 的距离的倒数来加权的。在预测时，距离越近的最近邻样本对最终决策的影响越大，距离越远的样本影响越小。

权重函数为K-近邻算法提供了一种灵活的决策方式。在某些情况下，我们可能希望距离较近的邻居对分类结果有更大的影响，这时可以选择距离权重函数。而在其他情况下，如果希望所有邻居对分类结果的影响相同，则可以选择统一函数。权重函数的选择可以根据具体问题的需求和数据集的特点来决定。

7.2.3 K-近邻算法应用与评价

K-近邻是一种懒惰学习算法，没有显式的学习过程，该算法不在训练时建立模型，而是在预测时才去查找最近邻居来做预测。此外，K-近邻是一种非参数化方法，不需要假设数据的分布形式，适应各种类型的数据集。然而，当训练集中的类别不平衡时，即某些类别的样本数量远远超过其他类别时，K-近邻可能会倾向于那些具有更多样本的类别，从而影响分类结果的准确性。

懒惰学习（Lazy learning）是机器学习中的一种策略，也称为基于实例的学习（Instance-based learning）或记忆型学习（Memory-based learning）。与传统的急切学习（Eager learning）算法相对应，懒惰学习的主要思想是将训练数据集存储在内存中，并在需要进行预测时直接使用已有的训练样本，而不是显式地构建模型，即根据训练数据来学习一个函数或模型，该模型能够描述输入特征和输出标签之间的关系。

K-近邻算法的优点在于实现简单、适应性强，但同时也存在一些局限性。由于K-近邻算法在预测时需要计算待预测样本与所有训练样本之间的距离，因此在大规模数据集上，计算量可能会非常大，导致预测速度较慢。此外，当数据集中的类别不平衡时，K-近邻算法可能会受到样本数量较多的类别的影响，从而影响分类结果的准确性。因此，在实际应用中，我们需要根据数据集的特点和问题的需求来合理选择和调整算法的参数，以获得更好的分类性能。

7.2.4 实践案例：基于K-近邻的广告点击预测

本节利用K-近邻算法对某网站投放的广告是否被用户点击进行分类预测。数据为数据科学家Jose Portilla和数据科学培训机构Pierian Data创建的广告点击数据集，该数据集包含1000条广告点击数据，具体字段包括：每日网站浏览时长、用户年龄、地区平均收入、每天上网时间、广告标题、城市、是否男性、国家、时间戳、是否点击广告。

在学习这个实践案例时，我体会到了数据预处理的重要性。为了避免量纲对距离计算的影响，在使用K-近邻算法前，需要对特征进行归一化。同时，对于离散特征如广告标题、城市、国家等，需要通过独热编码等方式将它们转换为数值特征。由于K-近邻是一种基于距离度量的分类算法，为了避免编码给样本距离度量带来的误差，本例在数据预处理阶段剔除了广告标题、城市、国家、时间戳这四个特征。对于保留的特征，使用scikit-learn中的StandardScaler类实现对特征的z-score归一化。

7.3 朴素贝叶斯分类

朴素贝叶斯分类是一种基于贝叶斯定理的统计学分类方法，它在许多领域有着广泛应用，例如文本分类、垃圾邮件过滤和医学诊断等。通过学习朴素贝叶斯分类，我对其基本原理、算法实现以及应用有了更深入的理解。

7.3.1 贝叶斯分类基本原理

贝叶斯分类的核心在于利用贝叶斯定理来计算后验概率，从而确定样本的类别。贝叶斯定理描述了两个条件概率之间的关系，通过已知的先验概率和新的证据（条件概率），可以计算出更新后的后验概率。

- **先验概率 (Prior probability)**：在不考虑任何新证据之前，根据以往的知识或经验，对事件概率的初始估计。例如，在邮件分类中，假设历史信息表明100封邮件中有20封邮件为垃圾邮件，那么先验概率 $P(\text{垃圾邮件}) = 20/100 = 0.2$ 。
- **条件概率 (Conditional probability)**：在给定某个条件下，事件发生的概率。以邮件分类为例，在已知某邮件为垃圾邮件的情况下，该邮件中存在“在家挣钱”关键词的概率，即 $P(\text{在家挣钱} | \text{垃圾邮件})$ ，就是条件概率。
- **后验概率 (Posterior probability)**：在考虑了新的信息或证据之后，事件发生的概率。在已知某邮件中存在“在家挣钱”关键词的情况下，该邮件为垃圾邮件的概率，即 $P(\text{垃圾邮件} | \text{在家挣钱})$ ，就是后验概率。

贝叶斯定理给出了计算后验概率的方法： $P(c|x) = \frac{P(x|c) \times P(c)}{P(x)}$ ，其中 $P(c)$ 是先验概率， $P(x|c)$ 是条件概率， $P(c|x)$ 是后验概率。在分类问题中，我们需要确定 $P(c|x)$ ，即给定观测数据样本的特征 x ，该样本类别为 c 的概率。

先验概率和条件概率的确定是关键。先验概率通常可以通过历史数据统计得到，而条件概率的计算则需要根据具体问题和数据集的特点来进行。贝叶斯定理为我们提供了一个框架，通过结合先验知识和新的证据，能够对事件的概率进行更新和预测。

7.3.2 朴素贝叶斯分类原理

朴素贝叶斯分类对特征条件概率分布 $P(x|c)$ 做了条件独立性假设，即假设样本的各个特征之间相互独立。在这一假设下，联合概率 $P(x|c)$ 可以分解为各个特征条件概率的乘积： $P(x|c) = \prod_{i=1}^m P(x_i|c)$ ，其中 m 是样本的特征数量， x_i 是 x 的第 i 个特征的值。

朴素贝叶斯分类的预测准则为：选择使后验概率 $P(c|x)$ 最大的类别 c 作为预测结果。由于 $P(x)$ 对所有类别相同，因此可以简化为： $\hat{c} = \arg\max_c P(c) \times \prod_{i=1}^m P(x_i|c)$ 。

在实际分类任务中，先验概率 $P(c)$ 和条件概率 $P(x_i|c)$ 均可以从训练集中估计得到。对于离散特征，条件概率可以通过统计训练集中特征值的频数来计算；对于连续特征，通常假设其服从高斯分布，从而计算条件概率。

条件独立性假设是其核心，它大大简化了概率计算的复杂度。然而，这一假设在现实世界中往往不成立，因为特征之间可能存在相关性。尽管如此，朴素贝叶斯分类在许多情况下仍然表现出色，尤其是在特征数量较多且样本数量较少的情况下。这可能是因为大量特征的情况下，即使存在一些相关性，独立性假设带来的误差也会被稀释。

7.3.3 零频现象的拉普拉斯修正

在朴素贝叶斯分类中，当某个特征值在训练集中没有与某个类同时出现过时，会出现“零频现象”，导致条件概率为零。例如，在员工离职数据中，不存在“考核=差，离职=否”的样本，因此 $P(\text{考核=差}|\text{离职=否})=0$ 。这会导致基于该特征的预测结果不准确。

为了解决零频现象，通常使用拉普拉斯修正。拉普拉斯修正通过给所有未出现的特征值赋予一个小的非零权重来避免概率估计中的零概率情况。具体来说，对于离散特征，条件概率的计算公式修正为：
$$P(x_i|c) = \frac{|D_{c,i}| + 1}{|D_c| + N_i}$$
，其中 $|D_{c,i}|$ 是训练集中类别为 c 且特征 x_i 取特定值的样本数量， $|D_c|$ 是训练集中类别为 c 的样本数量， N_i 是特征 x_i 可能的取值数。

拉普拉斯修正能够有效平衡已观察到的事件和未观察到的事件之间的概率分配，提供更平滑的估计结果。它在处理稀疏数据和小样本数据时尤为重要，能够避免因数据不足导致的预测偏差。

拉普拉斯修正能够在不引入过多复杂度的情况下，解决概率估计中的零频问题。这让我认识到在实际应用中，对于数据的处理和模型的优化，往往需要一些巧妙的技巧和方法，而不仅仅是依赖于复杂的算法。

7.3.4 朴素贝叶斯算法应用与评价

朴素贝叶斯算法基于概率原理，实现相对简单，适合处理大规模数据集。此外，它对缺失数据不太敏感，即使某些特征存在数据缺失，也能够进行分类。然而，朴素贝叶斯算法有严格的独立性假设，以及对数据分布敏感，这些因素可能导致算法性能下降且无法学习特征之间的相互作用。

在实际应用中，朴素贝叶斯算法的优点在于其简单性和高效性。它不需要复杂的模型构建过程，只需计算先验概率和条件概率，即可进行分类预测。这使得它在处理大规模数据集时具有较高的效率。同时，朴素贝叶斯算法对缺失数据的容忍度较高，能够处理不完整的数据集，这在实际数据收集过程中具有一定的优势。

然而，朴素贝叶斯算法的局限性也不容忽视。由于其条件独立性假设，当特征之间存在较强的相关性时，算法的性能可能会受到影响。此外，对于数据分布的假设（如高斯分布）也可能不完全符合实际情况，从而影响分类结果的准确性。因此，在实际应用时，需要根据数据集的特点和问题的需求，合理选择和调整算法的参数，或者与其他算法结合使用，以获得更好的分类性能。

7.3.5 实践案例：基于朴素贝叶斯预测恒星类型

本节利用朴素贝叶斯算法对恒星进行分类预测，使用的是Kaggle的星型分类数据集。该数据集包含240个样本，具体字段包括：温度、相对光度、相对半径、绝对幅度、颜色、SMASS规格、类型。数据预处理阶段，需要对离散特征进行哑编码，并统一颜色特征的不同取值。模型构建时，使用scikit-learn中的高斯朴素贝叶斯分类器（GaussianNB）。

在实践案例中，我体会到了数据预处理的重要性。对于离散特征，通过哑编码将其转换为数值型特征，使得模型能够更好地处理这些数据。同时，对于颜色特征的不同取值进行统一处理，避免了因数据不一致导致的模型训练问题。此外，通过可视化数据样本，可以更直观地了解数据的分布情况，为模型的选择和优化提供参考。

模型训练和测试结果显示，使用高斯朴素贝叶斯算法对恒星分类的准确率为0.933，表明该算法在该数据集上具有较好的分类性能。然而，仍有4个样本被预测错误，这提示我们在实际应用中，还需要进一步优化模型，提高分类的准确性。例如，可以尝试调整模型的参数，或者结合其他特征选择和降维方法，以更好地挖掘数据中的有用信息，提高模型的预测能力。

朴素贝叶斯算法以其简单、高效的特点，在许多实际问题中得到了广泛应用。然而，它也存在一些局限性，需要我们在实际应用中结合具体情况进行分析和优化。未来可以进一步探索其他概率分类方法，以及如何将朴素贝叶斯算法与其他算法相结合，以解决更复杂的数据分类问题。

7.4 决策树

决策树是一种基于树形结构的分类与回归算法，它模拟了人类在面对决策时的思维过程，帮助人们理清各个条件的权重和影响，从而更好地做出决策。

7.4.1 决策树基本原理

7.4.1.1 结构组成

决策树由一个根结点、若干内部结点、若干叶结点和若干有向边构成。根结点代表整个数据集，内部结点表示特征为某特定值时的数据子集，叶结点对应最终的分类结果。通过一系列特征选择和判定，决策树逐步分割数据集形成树形结构。在建立完成的决策树中，待预测样本从根结点开始，经过内部结点判断特征，最终到达叶结点得出分类结果。

在学习决策树的结构时，我意识到它与人类的决策过程非常相似。我们常常会根据一系列条件来进行判断和选择，而决策树将这一过程形式化和可视化。这种结构使得决策树易于理解和解释，我们可以清晰地看到每个特征在决策过程中的作用和影响。

决策树

Decision Tree

决策树是一种基于树结构的算法，用于分类和回归任务。其工作原理类似于人类做决策的方式——逐步根据条件分支，直到得出最终答案

按照条件分裂方式可以分为：

二叉决策树



每个内部节点只能有两个子节点

多叉决策树

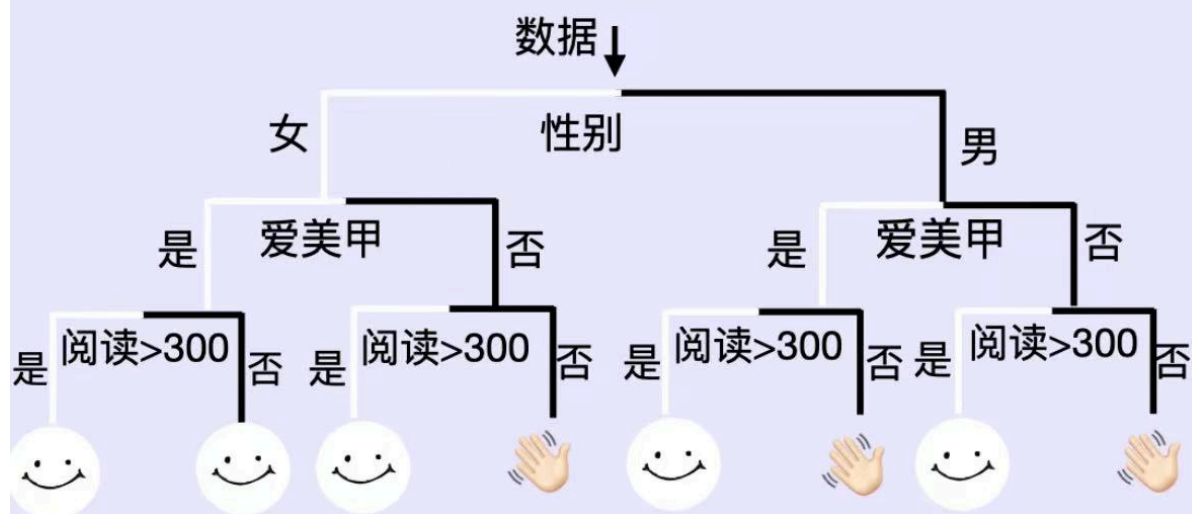


每个内部节点可以有多个子节点

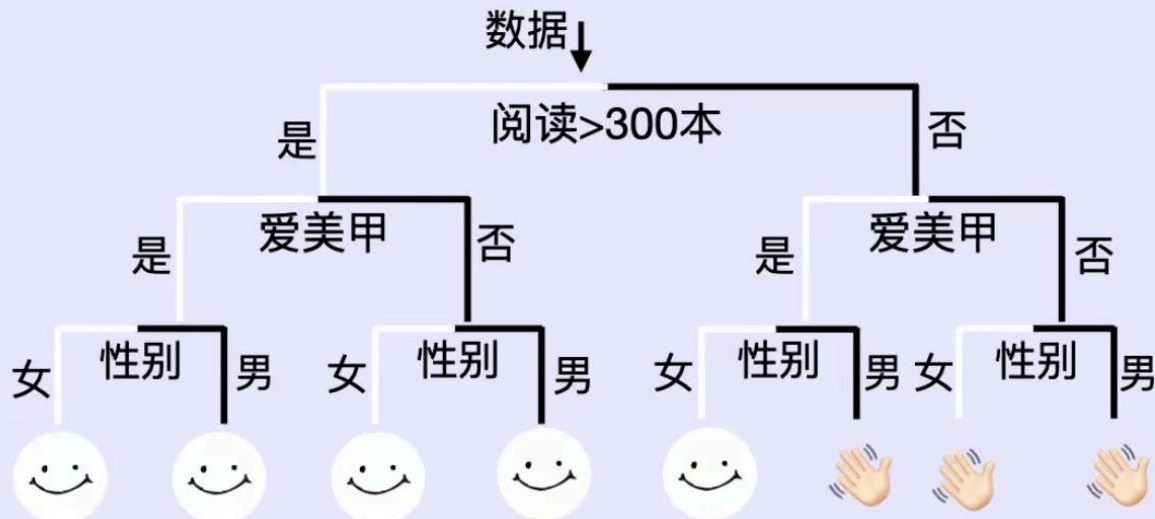
决策树

Decision Tree

二叉决策树：节点分裂只能二选一，节点分为是或否，
阈值大于或小于，性别男和女等二元分类



更换切分顺序为：阅读>300本 -> 爱美甲 -> 性别，则：



决策树

Decision Tree

多叉决策树：灵活，支持不同类别



特点	二叉决策树	多叉决策树
分裂方式	只能分成两个子节点	可分成多个子节点
计算复杂度	较小	较大
适用情况	连续变量和二分类任务	多类别分类任务
决策效率	每次决策只需判断是/否，速度快	需要同时评估多个选项，速度较慢
易解释性	结构清晰	较复杂，路径更多
过拟合风险	低（如果树深度适中）	较高（如果分支太多，可能导致数据过拟合）
决策树路径	长，要进行多个是/否判断	短，多类别可以直接分开，更高效
剪枝优化	剪枝避免树过深	剪枝避免树过宽

7.4.1.2 特征选择

特征选择是决策树构建的关键步骤之一。通过特征的选择和分裂，可以将数据集分割为不同的子集，使得子集的纯度提高（即不确定性降低）。度量数据集纯度的方法有多种，如信息增益、信息增益率和基尼系数等。

- **信息增益**：衡量一个特征对数据集划分后，信息复杂度（不确定性）的减少程度。信息增益越大，说明根据该特征对数据进行划分后的数据子集纯度越高，该特征越适合作为决策树中的分裂结点。
- **信息增益率**：对信息增益进行修正，以解决信息增益倾向于选择具有较多取值的特征的问题。信息增益率考虑了特征的固有值，能够更准确地评估特征的划分能力。
- **基尼系数**：度量一个数据集中随机选取两个不同样本时，这两个样本在类别上不一致的概率。基尼系数越小，表示数据集的纯度越高。

信息增益简单直观，但存在偏置问题；信息增益率能够克服这一问题，但计算较为复杂；基尼系数计算简单，适用于CART算法。选择合适的特征选择方法需要根据具体问题和数据集的特点来决定。

7.4.1.3 决策树生成

决策树的生成是一个递归地选择最优特征，并根据该特征对训练数据进行分割的过程。初始阶段，根据所有训练数据构建根结点。然后，选择最优特征将训练数据集分割成子集。如果子集中的样本能够基本正确分类（纯度高于指定阈值），则构建叶结点，并将子集分到相应的叶结点中；否则，选择新的最优特征，继续分割并构建相应的结点。这一过程将递归进行，直至达到终止条件之一。

通过不断地分割数据集，决策树能够逐步细化分类规则，直至找到合适的分类依据。然而，这一过程也容易导致过拟合，因此需要引入剪枝技术来优化决策树。

7.4.1.4 剪枝

决策树在训练过程中容易过拟合，剪枝是一种用于减少模型复杂度的技术，具体可分为预剪枝和后剪枝两种类型。

- **预剪枝 (Pre-pruning)**：在构建决策树的过程中，在每个结点处判断是否停止分裂，以避免过拟合。通过设置条件，如限制最小样本数或最大深度，阻止决策树过度划分。
- **后剪枝 (Post-pruning)**：在树已经构建完成后，通过去除一些结点或分支来实现剪枝。通常从叶结点开始逐步向根结点的方向剪枝，剪枝时使用测试集评估剪枝后的决策树性能，若表现良好则保留剪枝后的决策树，否则保留原始决策树。

剪枝对于提高决策树的泛化能力至关重要。适当的剪枝能够去除不必要的分支，简化决策树结构，使其更加简洁和易于理解。同时，剪枝也有助于提高决策树在新数据上的预测性能，避免因过拟合而导致的模型失效。

7.4.2 ID3 算法

ID3算法的核心思想是通过信息增益对特征进行选择。在ID3决策树的生成过程中会计算每个特征的信息增益，选择给数据集带来的信息增益最大的特征作为决策树当前划分的结点。需要注意的是，ID3决策树不存在剪枝操作。

7.4.2.1 信息增益

- **信息熵 (Information entropy)**：用于衡量一个随机变量的不确定性。假设训练数据集D中有K类不同的样本，C是类别为k的数据子集，则D的信息熵为： $H(D) = -\sum_{k=1}^K \frac{|C_k|}{|D|} \log_2 \frac{|C_k|}{|D|}$ ，其中，|D|表示数据集D的样本数，|C_k|表示类别为k的数据子集的样本数。
- **条件熵 (Conditional entropy)**：在给定另一个随机变量的条件下，一个随机变量的不确定性。假设决策树中的某个结点有一个特征A，该特征有d个不同的取值，根据特征A的不同取值将数据集划分为d个子集，则在给定特征A条件下D的条件熵为： $H(D|A) = \sum_{i=1}^d \frac{|D_i|}{|D|} H(D_i)$ ，其中，D_i表示第i个数据子集中类别为k的样本集合。
- **信息增益 (Information gain)**：衡量一个特征对数据集划分后，信息复杂度（不确定性）的减少程度的指标。对于一个数据集D和特征A，其信息增益为： $\text{Gain}(D|A) = H(D) - H(D|A)$ 。

通过计算信息增益，ID3算法能够选择最优特征进行数据划分，从而构建出一棵决策树。然而，信息增益存在一个内生性的偏置，它更倾向于选择具有较多取值的特征作为划分依据，这可能导致模型在未来的预测任务中过拟合。

7.4.3 C4.5 算法

C4.5算法是从ID3算法演变而来，该算法使用信息增益率作为特征选择的依据来克服信息增益存在的问题。

7.4.3.1 信息增益率

信息增益率 (Gain ratio) 特征A对训练数据集D的信息增益比 $\text{GainRatio}(D|A)$ 为信息增益 $\text{Gain}(D|A)$ 与D关于特征A的熵 $H(A)$ 之比，其计算为：
$$\text{GainRatio}(D|A) = \frac{\text{Gain}(D|A)}{H(A)}$$
，其中， $H(A) = -\sum_{i=1}^d \frac{|D_i|}{|D|} \log_2 \frac{|D_i|}{|D|}$ ，d为特征A取值的个数。

C4.5决策树使用信息增益率作为特征划分的依据，其生成的具体流程与ID3类似，但在特征选择时采用信息增益率代替信息增益。此外，C4.5算法还引入了剪枝操作，使用悲观错误剪枝 (Pessimistic Error Pruning) 的方法对决策树进行剪枝。

7.4.3.2 悲观错误剪枝

悲观错误剪枝的基本思想是当决策树的某个分支被剪掉时，模型在未来数据上的性能可能会下降。因此，悲观误差剪枝会计算不剪枝情况下模型在未来数据上的性能的一个悲观估计，并与剪枝模型的性能进行比较。如果剪枝后的误差不大于最悲观误差，则执行剪枝操作。

信息增益率的引入有效地解决了信息增益的偏置问题，使得特征选择更加合理。同时，悲观错误剪枝为决策树的优化提供了有力的手段，能够在一定程度上提高模型的泛化能力。

7.4.4 CART 算法

CART (Classification and Regression Trees) 决策树又称分类回归树，它既可用于分类，又可用于回归。在分类任务中，CART决策树使用基尼系数衡量数据集的纯度。

7.4.4.1 基尼系数

基尼系数 (Gini index) 是用于度量一个数据集中随机选取两个不同样本时，这两个样本在类别上不一致的概率。

7.4.4.2 二叉分裂

CART生成的决策树是一棵二叉树，每个结点都是一个二元判定规则，将输入样本分成两部分。在决策树的生成方面，CART通过计算基尼系数来选择最优特征及其对应的切分点，从而进行二叉分裂。

7.4.4.3 代价复杂度剪枝

CART采用基于代价复杂度剪枝 (Cost complexity pruning) 的方法对决策树进行剪枝。代价复杂度剪枝的损失函数包含树T的代价 (Cost) 和复杂度，以及衡量代价和复杂度之间关系的参数 $\alpha \geq 0$ 。对于决策树T，假设t为树T的叶结点，该结点对应的样本数量为 N_t ，其中第k类样本数量为 N_{tk} ， $k=1,2,\dots,K$ 。 $E_t(T)$ 为叶结点对于样本的信息熵， $\alpha \geq 0$ 为参数。在参数为 α 下，剪枝的损失函数 $C_\alpha(T)$ 为：
$$C_\alpha(T) = C(T) + \alpha |T|$$
，其中， $C(T)$ 表示模型对训练数据的预测误差， $|T|$ 表示模型的复杂度 (可用叶结点数量表示)，参数 α 则用于控制二者的影响程度， α 越大，模型越简单，反之，模型越复杂。

基尼系数和二叉分裂是其核心概念。基尼系数为特征选择提供了有效的衡量标准，而二叉分裂使得决策树的结构更加简洁。代价复杂度剪枝则在模型的复杂度和预测性能之间找到了一个平衡点，有助于提高决策树的泛化能力。

7.4.5 决策树算法应用与评价

ID3、C4.5和CART是三种常用的决策树算法，它们在特征选择、决策树生成和剪枝上有着不同的准则。

- **ID3算法**：易于理解和实现，生成的决策树通常具有清晰的结构。然而，没有剪枝操作，可能导致生成的决策树过于复杂，泛化能力较差。同时，倾向于选择具有更多取值的特征作为划分依据，这也将导致其生成的决策树泛化能力较差。
- **C4.5算法**：使用信息增益率作为特征划分标准，克服了信息增益存在的内在偏置问题。引入剪枝操作，故在时间开销上略低于ID3。然而，信息增益率需要进行极为耗时的对数运算。
- **CART算法**：采用二叉树结构，选用基尼系数或均方差等指标进行特征选择，训练时间相对较快。然而，生成的树是二叉树，每个结点只能有两个子结点，这可能不足以捕捉某些复杂的关系。

在实际应用中，选择合适的决策树算法需要根据数据集的特点和问题的需求来决定。例如，对于数据集较小且特征数量较少的问题，ID3算法可能是一个不错的选择；而对于需要处理大规模数据集且对模型泛化能力要求较高的问题，C4.5或CART算法可能更为合适。此外，还可以通过调整算法的参数或与其他算法结合使用，以获得更好的分类性能。

7.4.6 实践案例：基于决策树的旅游业客户流失预测

本节利用CART决策树算法对旅游业客户流失进行预测，数据集来源于一家旅游公司提供的954名客户信息。数据预处理阶段，对离散特征进行编码，并使用众数填充法对缺失值进行处理。模型训练时，使用scikit-learn中的DecisionTreeClassifier类构建CART决策树模型，并通过网格搜索寻找最佳参数组合。模型评价时，考虑到数据集的不平衡性，选择召回率作为评价指标。

在实践案例中，我体会到了数据预处理的重要性。对于离散特征，通过编码将其转换为数值型特征，使得模型能够更好地处理这些数据。同时，对于缺失值的处理也至关重要，合理的缺失值填充方法能够提高模型的训练效果。此外，通过网格搜索寻找最佳参数组合，可以有效地优化模型的性能。最终，模型在测试集上的召回率为0.815，表明该决策树模型能够较好地识别出可能流失的客户，为旅游业客户流失预测提供了一种有效的解决方案。

决策树以其直观、易于理解的特点，在许多实际问题中得到了广泛应用。然而，它也存在一些局限性，如容易过拟合、对数据中的噪声敏感等。因此，在实际应用中，需要根据具体情况进行分析和优化，以充分发挥决策树的优势。

7.5 Logistic回归

Logistic回归是一种经典的分类方法，它在数据分析和机器学习领域有着广泛的应用。通过学习Logistic回归，我对分类问题的处理有了更深入的理解。

7.5.1 从线性回归到Logistic回归

7.5.1.1 Logistic函数

Logistic函数是Logistic回归的核心，其图像呈现为一个S形曲线，定义域为 $(-\infty, +\infty)$ ，值域为 $(0, 1)$ 。当 $x > 0$ 时， $\sigma(x) > 0.5$ ；当 $x < 0$ 时， $\sigma(x) < 0.5$ 。Logistic函数的数学表达式为：
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

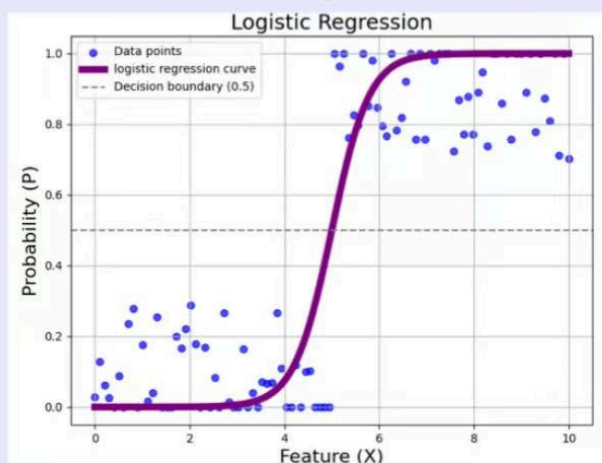
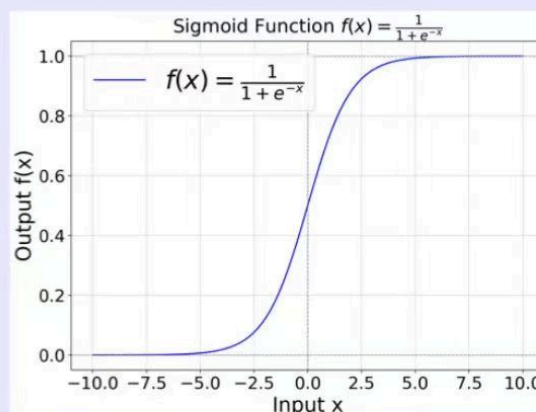
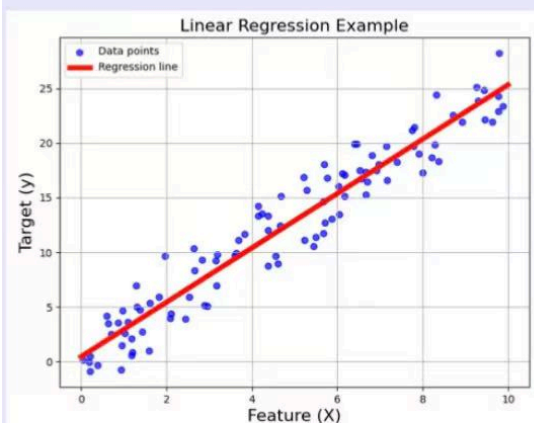
Logistic函数为线性回归输出的概率值提供了一个合适的映射。通过Logistic函数，我们可以将线性回归的结果转换为介于0和1之间的概率值，从而实现二分类问题的预测。这种转换使得Logistic回归能够处理分类问题，而不仅仅是回归问题。

逻辑回归

Logistic Regression

逻辑回归(Logistic Regression), 机器学习常见算法,
尽管名字里带“回归”, 但主要用于二分类问题。

线性回归 + sigmoid => 逻辑回归



逻辑回归

Logistic Regression

数学公式：

Sigmoid

线性回归

$$P(y = 1 | x) = \sigma(\underline{w^T x + b})$$

$$P(y = 0 | x) = 1 - \sigma(\underline{w^T x + b})$$

其中：

$P(y = 1 | x)$ ：样本属于正类的概率，

$P(y = 0 | x)$ ：样本属于负类的概率

x ：输入特征向量， w ：回归系数（权重）， b ：偏置项

$\sigma(z) = \frac{1}{1 + e^{-z}}$ ：逻辑函数（Sigmoid 函数），将线性

结果 $z = w^T x + b$ 转化为概率

根据概率值，逻辑回归通过设定阈值（如 0.5）将输入样本划分为正类或负类：

•如果 $P(y = 1 | x) \geq 0.5$ ，预测 $y = 1$ 。

•如果 $P(y = 1 | x) < 0.5$ ，预测 $y = 0$ 。

逻辑回归

Logistic Regression

损失函数

逻辑回归使用**交叉熵损失函数 (Cross-Entropy Loss)** 来衡量预测概率与真实标签的差异：

$$\mathcal{L}(y, \hat{y}) = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

其中：

y_i : 真实标签 (0 或 1)

\hat{y}_i : 预测概率 ($P(y = 1 | x)$)

N : 样本数。

目标是通过优化算法（如梯度下降）最小化损失函数，调整参数 w 和 b 。

(注：线性回归使用 MSE 是因为它预测的是连续值，MSE 对连续值的误差衡量效果很好。而逻辑回归输出的是概率，用交叉熵更合适，因为它能更好地捕捉预测概率与真实标签之间的差异，同时提供更有用的梯度用于优化)

逻辑回归

Logistic Regression

优点：

- 1.简单易用：**公式简单，容易实现。
- 2.高效：**计算复杂度低，适合大规模数据。
- 3.可解释性强：**回归系数反映特征对分类结果的影响
- 4.稳健性好：**适合二分类问题，效果稳定。

缺点：

- 1.假设线性关系：**只能处理线性可分数据，对复杂数据的表现有限。
- 2.对异常值敏感：**异常值可能对模型参数产生较大影响
- 3.特征工程依赖强：**需要对特征进行认真选择和处理

7.5.1.2 事件优势比

事件的优势比 (Odds ratio) 是Logistic回归中一个重要的概念，它表示事件发生的概率与事件不发生的概率的比值。对事件的优势比取自然对数即可得到Logistic变换：
$$\ln(\text{Odds}) = \ln\left(\frac{p}{1-p}\right)$$
，其中 $p=P(y=1|X)$ 。

事件优势比的可解释性在于，在Logistic回归中，每个特征的事件优势比表示在其他特征保持不变的情况下，该特征变化对事件发生几率的影响。若事件优势比为1，表示该特征对事件发生无影响；若大于1，则表示正向影响，即随该特征增加，事件几率增加；若小于1，则表示负向影响，即随该特征增加，事件几率减小。

事件优势比为我们提供了一种直观的方式来理解各个特征对分类结果的影响。通过分析事件优势比，我们可以判断哪些特征对分类结果具有显著影响，从而为特征选择和模型优化提供依据。

7.5.1.3 Logistic回归模型构建

Logistic回归模型是建立在线性回归模型的基础上，通过Logistic函数将线性回归的输出映射到(0,1)的概率空间。其基本形式为： $p = \sigma(\theta^T x)$ ，其中 x 是样本的特征向量， θ 是权重向量， σ 是Logistic函数。

在模型构建过程中，我们需要确定合适的权重向量 θ ，使得模型能够准确地预测样本的类别。这通常通过最小化损失函数来实现。

7.5.2 Logistic回归的损失函数

Logistic回归模型的求解目标是找到合适的参数 θ ，使得模型能够最大程度地对样本进行正确分类。为此，我们需要构建Logistic回归的损失函数。

给定 n 个样本的特征 X 和对应类别 Y ，Logistic回归的损失函数为： $J(\theta) = -\frac{1}{n} \sum_{i=1}^n [y_i \ln(p_i) + (1-y_i) \ln(1-p_i)]$ ，其中 p_i 是样本 x_i 为正样本的概率。

损失函数的构建基于最大似然估计的思想，通过最大化样本的似然函数来确定模型参数。在Logistic回归中，损失函数的最小化等价于似然函数的最大化。

通过最小化损失函数，我们可以找到最优的模型参数，使得模型在训练数据上的预测误差最小。然而，Logistic回归的损失函数并没有显式的解析解，因此需要采用数值优化方法来求解。

7.5.3 通过梯度下降求解最优参数

由于Logistic回归的损失函数没有显式的解析解，通常使用梯度下降法求解模型参数。梯度下降是一种用于寻找函数最小值的优化算法，其基本步骤如下：

1. 随机初始化模型参数 θ 。
2. 计算损失函数关于参数 θ 的梯度，即损失函数在当前参数值最快减小的方向。
3. 更新参数 θ ， $\theta = \theta - \alpha \nabla J(\theta)$ ，其中 α 是学习率，用于控制每次参数更新的步长。
4. 重复执行步骤2和步骤3，直到达到最大迭代次数或达到预定的精度水平。

梯度下降法是一种简单而有效的优化方法。通过不断迭代更新参数，梯度下降法能够逐渐逼近损失函数的最小值，从而找到最优的模型参数。然而，梯度下降法的收敛速度和最终结果受到初始参数选择和学习率设置的影响，因此在实际应用中需要进行适当的调整和优化。

7.5.4 OvR 和 OvO

Logistic回归主要用于二分类问题，但在实际应用中，我们常常需要对多个类别进行分类。为此，可以采用一些特定的策略使二元分类模型能够完成多分类任务，常用的方法包括OvR (One-vs-Rest) 和OvO (One-vs-One)。

- **OvR**：对于具有 N 个类别的问题，将其转化为 N 个二元分类问题。对于每个类别，训练一个分类器来区分该类别与其他所有类别的组合。在预测阶段，使用 N 个分类器分别进行预测，将置信度最高的类别作为最终的预测结果。
- **OvO**：对于具有 N 个类别的问题，将其转化为 C^2 个二元分类问题，即针对“每对”类别训练一个分类器。在预测阶段，通过对所有分类器的结果进行投票或应用其他决策规则来确定最终的预测结果。

OvR和OvO为Logistic回归在多分类问题中的应用提供了有效的解决方案。OvR方法简单易实现，适用于类别较多的情况；而OvO方法在类别较少时可能更为准确。具体选择哪种方法需要根据实际问题的数据规模和类别特点来决定。

7.5.5 Logistic回归算法应用与评价

Logistic回归算法具有简单易实现、计算效率高、可解释性强、适用于稀疏数据和不容易过拟合等优点。然而，它也存在一些局限性，包括对非线性关系的建模能力相对较弱、对多重共线性敏感等。在处理复杂的非线性问题或存在高度相关性的特征时，Logistic回归的表现可能不如其他更复杂的模型。

在实际应用中，Logistic回归常用于医疗诊断、信用评估、市场营销等领域。例如，在医疗诊断中，Logistic回归可以用于预测患者是否患有某种疾病；在信用评估中，它可以用于评估客户的信用风险。尽管Logistic回归在某些情况下可能不如其他模型表现优异，但其简单性和可解释性使其在许多实际问题中仍然具有重要的应用价值。

7.5.6 实践案例：基于Logistic回归的肝病预测

本节利用Logistic回归算法对病人是否患有肝病进行预测，使用的数据集是来自印度安得拉邦东北部的416位肝病患者记录和167位非肝病患者的部分医学指标记录。

在实践案例中，我体会到了数据预处理的重要性。对于离散特征，通过哑编码将其转换为数值型特征，使得模型能够更好地处理这些数据。同时，对于缺失值的处理也至关重要，合理的缺失值填充方法能够提高模型的训练效果。此外，特征归一化有助于提高模型的收敛速度和预测性能。

模型训练时，使用scikit-learn中的LogisticRegression类构建Logistic回归模型，并设置合适的正则化参数C。正则化参数C的大小会影响模型的复杂度和过拟合程度，需要根据具体问题进行调整。在本例中，通过设置C=0.1，模型在测试集上的精准率为0.926，表明该Logistic回归模型能够较好地预测患者是否患有肝病。

通过分析事件优势比，我们可以了解每个特征对于结果的影响。例如，“直接胆红素”的事件优势比高达4.85，意味着理论上直接胆红素每增加一个单位，患肝病的几率将增加4.85倍。这种分析有助于我们识别出对肝病预测具有重要影响的特征，为临床诊断提供参考。

通过学习Logistic回归及其实例案例，我对分类问题的处理有了更深入的理解。Logistic回归以其简单、高效的特点，在许多实际问题中得到了广泛应用。然而，它也存在一些局限性，需要我们在实际应用中结合具体情况进行分析和优化。

7.6 支持向量机

支持向量机（Support Vector Machines, SVM）是一种强大的二分类模型，它通过在特征空间中找到一个最大间隔超平面来划分不同类别的数据，从而实现分类。此外，SVM还可以通过核函数将数据映射到高维空间以解决非线性分类问题。通过学习SVM，我对分类问题的处理有了更深入的理解。

7.6.1 支持向量机概述

SVM的主要思想

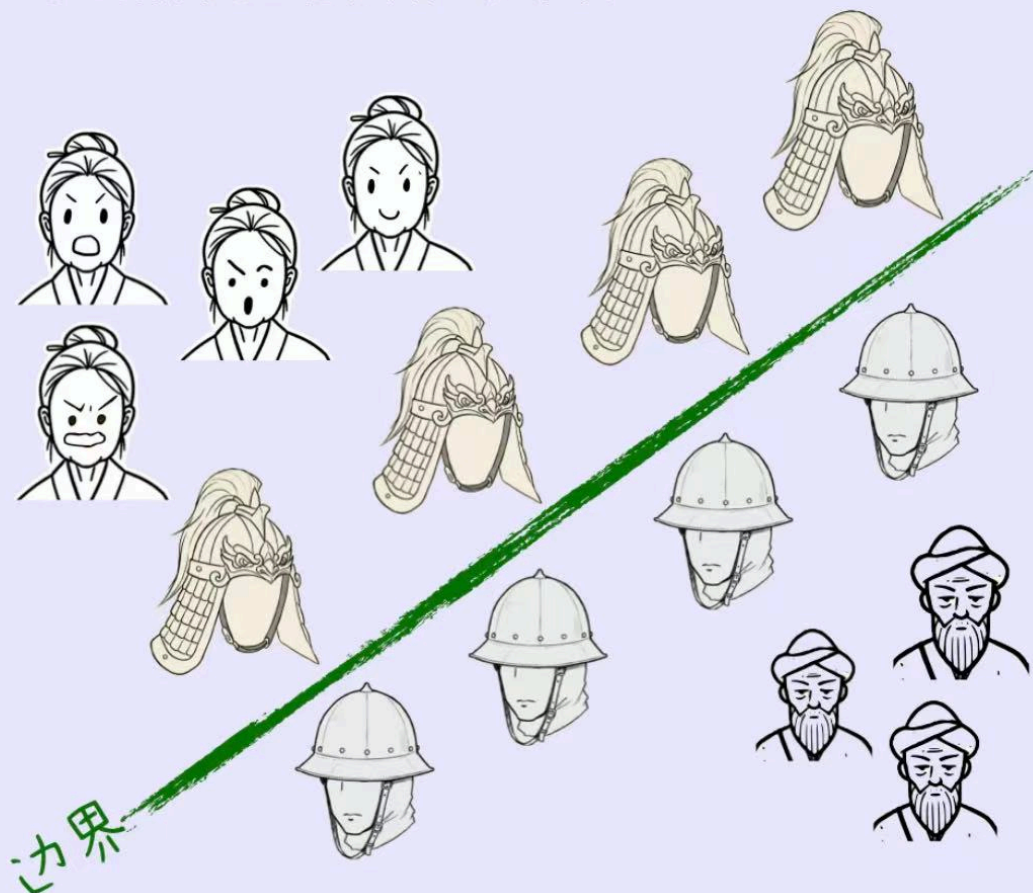
SVM的核心思想是寻找能够有效区分不同类别的决策边界（或超平面），并使得这个决策边界与各类别的最近数据点之间的间隔最大化。这种最大间隔的决策边界有助于提高模型的泛化能力，从而在新数据上表现更为优越。

间隔最大化是其独特的优势之一。与其他分类模型相比，SVM通过最大化间隔来提高模型的鲁棒性和泛化能力，这使得它在处理线性可分和非线性可分数据时都能取得良好的效果。

支持向量机

SVM: Support Vector Machine

支持向量机（**SVM: Support Vector Machine**），一种常见的机器学习算法，常被应用于二分类问题，同时也可以拓展到多分类或回归任务



但使龙城飞将在 不教胡马度阴山

支持向量机

SVM: Support Vector Machine

核心思想

1. 寻找最优决策边界（超平面）

对于二分类问题，SVM 的目标是寻找一个“超平面”将数据分成两类，并找到**最大间隔**，尽量让“最近的样本点”与这个超平面之间的距离最大化



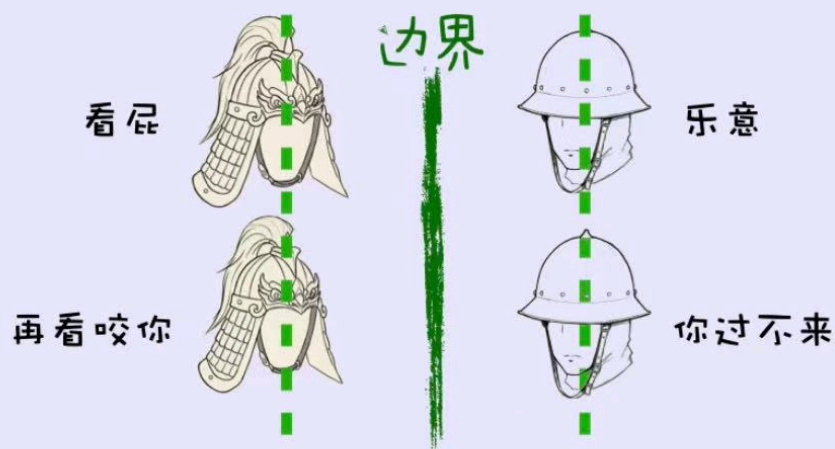
2. 间隔越大，模型泛化能力越好

最大化边界，SVM 能够减少在新数据上的错误分类概率

支持向量机

SVM: Support Vector Machine

与超平面“最近的样本点”就是 **支持向量**，参与决定最大间隔的运算（只有将士冲锋陷阵）



与线性回归，逻辑回归等算法不同的是，其他算法所有数据点都要参与运算，而SVM只有支持向量参与决策边界，非支持向量的点都不参与决策边界



老娘不是不想上
只是有娃还要养



边关悬而未决
不如回家种田

支持向量机

SVM: Support Vector Machine

核函数 (Kernel Trick)

当数据在原有的特征空间中无法被线性分类时，可以使用 **核函数 (Kernel)** 将数据映射到高维或无限维空间，使得在这个空间里数据可以被线性可分。

常见的核函数包括：

1. **线性核** (Linear Kernel)
2. **多项式核** (Polynomial Kernel)
3. **径向基函数核** (RBF Kernel)
4. **Sigmoid 核** (Sigmoid Kernel)

SVM的模型分类

SVM从简单到复杂涵盖了一系列模型，包括线性可分SVM、线性SVM以及非线性SVM。简单模型不仅奠定了复杂模型的基础，同时也可视为复杂模型的特殊情况。

- **线性可分SVM**：假设训练数据集是线性可分的，即存在一个超平面能够将正样本和负样本完全分开。其目标是找到一个最优的超平面，使得所有正样本和负样本到这个超平面的距离最大。
- **线性SVM**：不假设训练数据集是线性可分的，即使数据集线性不可分，线性SVM仍然可以通过引入松弛变量来容忍一些分类错误。
- **非线性SVM**：用于处理非线性可分数据集的一种扩展形式。通过引入核函数将数据映射到高维空间，从而使其在高维空间中变得线性可分。

不同类型的SVM适用于不同场景的数据。对于线性可分的数据，线性可分SVM是最合适的选择；而对于线性不可分的数据，线性SVM和非线性SVM则提供了有效的解决方案。选择合适的SVM模型需要根据数据的特性和问题的需求来决定。

7.6.2 线性可分支持向量机

7.6.2.1 间隔表示

- **支持向量**：训练数据集中距离超平面最近的样本点，其直接决定了超平面的位置和方向。
- **间隔**：两类支持向量到超平面的距离之和。

最好的超平面应该位于两类训练样本的“正中间”，该划分超平面对训练样本局部扰动的“容忍性”最好，所产生的分类结果理论上是最鲁棒的，对未见实例的泛化能力是最强的。

支持向量和间隔是线性可分SVM的关键概念。支持向量决定了超平面的位置，而间隔的大小则反映了模型的鲁棒性和泛化能力。通过最大化间隔，SVM能够找到最优的超平面，从而实现数据的准确分类。

7.6.2.2 目标函数

训练线性可分SVM的目标是最大化“间隔”，即使得两类样本的支持向量到超平面的距离最大。目标函数可表示为： $\min_{w,b} \frac{1}{2} \|w\|^2$ ，约束条件为： $y_i(w^T x_i + b) \geq 1, \quad i=1,2,\dots,n$ 。

最大化间隔等价于最小化 $\|w\|^2$ ，这是因为间隔与 $\|w\|$ 成反比。通过优化目标函数，SVM能够找到最优的权重向量 w 和偏置项 b ，从而确定最优的超平面。

7.6.2.3 求解目标函数的对偶问题

求解线性可分SVM的目标函数是一个凸二次规划问题，原问题和对偶问题满足对偶关系。为了降低计算复杂度，可通过求解原问题的对偶问题得到原问题的解。

拉格朗日乘子法是求解凸优化问题的常用方法。通过构造拉格朗日函数并求解对偶问题，我们可以得到最优的拉格朗日乘子，进而求得最优的权重向量 w 和偏置项 b 。对偶问题的求解使得SVM的优化过程更加高效和稳定。

7.6.2.4 基于SMO算法求解目标函数

SMO (Sequential Minimal Optimization) 是一种启发式算法，其基本思想是每次只优化一个参数，并固定其他参数，仅求当前优化参数的极值。通过多次迭代直至收敛，求得最优解。

SMO算法通过每次只优化两个拉格朗日乘子，大大简化了优化过程，使得SVM的训练更加高效。此外，SMO算法的实现也相对简单，易于理解和应用。

7.6.3 线性支持向量机

7.6.3.1 软间隔

为了解决线性不可分的情况，提出了“软间隔”的概念。软间隔是指在SVM中允许一些样本点不被正确分类，即在求解最优间隔时允许一些样本点不严格满足约束条件。线性SVM不仅可以更好地防止过拟合，同时更不易受数据集中异常值的影响。

引入松弛变量 ξ 是实现软间隔的关键。松弛变量允许一些样本点位于超平面的错误一侧或位于间隔边界附近，从而使得模型能够更好地适应线性不可分的数据。软间隔的引入提高了SVM的灵活性和鲁棒性。

7.6.3.2 求解目标函数的对偶问题

线性SVM的对偶问题与线性可分SVM类似，但需要考虑松弛变量的影响。通过对偶问题的求解，我们可以得到最优的拉格朗日乘子，进而求得最优的权重向量 w 和偏置项 b 。

对偶问题的求解仍然是SVM优化的核心。通过对偶问题的求解，我们可以得到模型的最优参数，从而实现数据的准确分类。同时，对偶问题的求解也使得SVM能够处理大规模数据集，提高了模型的实用性。

7.6.3.3 基于SMO算法求解目标函数

与线性可分SVM类似，线性SVM也可以通过SMO算法求解目标函数。SMO算法通过每次只优化两个拉格朗日乘子，使得线性SVM的训练过程更加高效。

SMO算法的高效性对于处理大规模数据集至关重要。通过SMO算法，我们可以快速地求得线性SVM的最优参数，从而实现数据的准确分类。此外，SMO算法的实现也相对简单，易于理解和应用。

7.6.4 非线性支持向量机

7.6.4.1 目标函数

在现实任务中，原始样本空间内可能并不存在一个能正确划分两类样本的超平面。非线性SVM通过引入核函数将数据映射到高维空间，从而使其在高维空间中变得线性可分。

核函数是实现非线性映射的关键。通过核函数，我们可以将数据从低维空间映射到高维空间，从而使得原本线性不可分的数据变得线性可分。核函数的选择对于非线性SVM的性能有着重要的影响。

7.6.4.2 常用的核函数

- **线性核**：适用于数据线性可分的情况，特征空间与输入空间重合。
- **多项式核**：通过升维使得原本线性不可分的数据线性可分，但计算复杂度较高。
- **高斯核（径向基函数）**：通过计算两个样本之间的高斯距离来衡量样本之间的相似度，将原始的特征空间映射到更高维的特征空间。
- **Sigmoid核**：通过对样本之间的内积应用Sigmoid函数来衡量样本之间的相似度，但不能处理多分类问题。

不同核函数适用于不同场景的数据。线性核适用于线性可分的数据，多项式核和高斯核适用于非线性可分的数据，而Sigmoid核则适用于某些特定的场景。选择合适的核函数需要根据数据的特性和问题的需求来决定。

7.6.5 支持向量机算法应用与评价

SVM通过最大化间隔实现强泛化能力，在高维空间下有出色的性能表现，尤其适用于小样本和高维问题。在SVM中使用核函数使其能够处理非线性数据，同时可通过OvO或者OvR策略解决多类别问题。然而，SVM的计算复杂度较高，尤其在大规模数据集上训练时，模型参数调优相对困难，且对噪声较为敏感。

SVM在许多实际问题中得到了广泛应用，如图像识别、文本分类、生物信息学等。SVM的优势在于其强大的分类能力和良好的泛化性能，但同时也存在一些局限性。在实际应用中，需要根据具体问题的特点和数据集的规模来选择合适的模型，并进行适当的参数调优，以获得最佳的分类效果。

7.6.6 实践案例：不良用户识别

本节利用SVM算法对不良用户进行判断，数据集为Instagram上虚假账户和垃圾邮件发送账户的详细信息。数据预处理阶段，对特征进行归一化以消除特征数值范围对决策超平面的影响。模型训练时，使用scikit-learn中的SVC类构建SVM分类模型，并通过调整参数C、kernel、gamma等来优化模型性能。

在实践案例中，我体会到了数据预处理的重要性。特征归一化能够确保每个特征对模型的贡献是均衡的，从而提高模型的训练效果。同时，模型参数的选择对于SVM的性能有着重要的影响。通过调整参数C和gamma，我们可以控制模型的复杂度和过拟合程度，从而获得更好的分类结果。此外，核函数的选择也对模型的性能有着重要的影响，需要根据数据的特性和问题的需求来选择合适的核函数。

SVM以其强大的分类能力和良好的泛化性能，在许多实际问题中得到了广泛应用。然而，SVM也存在一些局限性，如计算复杂度较高、参数调优困难等。在实际应用中，需要根据具体问题的特点和数据集的规模来选择合适的模型，并进行适当的参数调优，以获得最佳的分类效果。

7.7 人工神经网络

人工神经网络（Artificial Neural Network, ANN）是一种受到生物神经网络启发而构建的计算模型。它由大量相互连接的神经元组成，这些神经元模拟了生物神经网络中的神经元之间的相互作用。这项技术已经在诸多领域中展现出惊人的成就，如图像识别、自然语言处理、智能控制等。

7.7.1 感知机

感知机是人工神经网络的基本单元，对应于生物神经网络中的神经元。它通过学习和调整连接权重，能够模拟信息的传递和处理过程，从而实现对复杂问题的学习和预测。

感知机模型结构

感知机由输入 $x=[x_1, x_2, \dots, x_n]$ 、偏置项 b 、输入权值 $w=[w_1, w_2, \dots, w_n]$ 、激活函数 f 和输出 y 组成。输入 x 表示感知机的 n 个输入，权值 w 表示神经元间的连接强度。每个输入值进行加权求和和运算的结果将输入到激活函数 f 中（在感知机中，激活函数为阶跃函数），激活函数的输出作为模型的最终输出 y 。

感知机模型结构是一种简单的线性二元分类器。通过线性组合输入特征并应用激活函数，感知机能够将特征向量映射到一个二元输出，即正类别或负类别。这种结构使得感知机在处理线性可分问题时具有一定的优势，但同时也限制了其在复杂问题中的应用。

感知机的学习规则

感知机模型的训练过程就是不断调整输入权值以获得期望的输出。感知机的学习规则是权值修正规则，具体步骤如下：

1. 随机初始化权重向量 w 和偏置项 b 。
2. 对于训练集 D 中的每一个样本 x_i ，计算感知机对该样本的预测值 y_i 。
3. 将感知机的预测值与真实值 t_i 进行比较，计算误差 $e_i = t_i - y_i$ 。
4. 根据 $w_i = w_i + \alpha e_i x_i$ 更新权重和偏置，使误差逐渐减小， α 表示学习率。
5. 重复步骤2到步骤4，直到达到预定的最大迭代次数 T 或误差小于阈值。

通过不断调整权重和偏置，感知机能够逐渐减小预测误差，从而实现对数据的准确分类。然而，感知机的学习能力有限，对于非线性可分问题，感知机无法找到合适的权重和偏置来实现准确分类。

7.7.2 XOR问题

感知机通过不断地调整权重 w ，可将线性可分的两类样本成功分隔开。因此，通过有限次的训练，感知机必定可以解决线性可分的问题。然而，对于异或问题，感知机却无法从输入数据中学习出适当的权重和偏置来解决这个问题。这是因为感知机只能产生线性的函数形式，而异或问题不是线性可分的。

异或问题是一个经典的非线性问题，它无法通过简单的线性函数来解决。这促使人们思考如何改进感知机模型，使其能够处理更复杂的非线性问题。

7.7.3 多层感知机模型

7.7.3.1 基本概念

为了克服感知机的局限性，多层感知机（Multilayer Perceptron, MLP）应运而生。多层感知机是一种人工神经网络的模型，它由多个神经元层组成，包括输入层、至少一个或多个隐藏层和输出层。输入层用于接收外部输入，通常由输入特征组成；隐藏层用于处理输入层提供的信息，将输入信息转换为更加抽象的特征，通常由多个神经元组成；输出层用于输出模型的预测结果。

多层感知机通过引入隐藏层，允许数据在多个非线性变换之间传递，从而捕捉到更为复杂的模式和关系。随着隐藏层数的增加，可以引入更多的非线性映射，使得模型可以解决更加复杂和抽象的问题。

隐藏层是其核心组成部分。隐藏层的存在使得多层感知机能够进行复杂的特征提取和表示学习，从而实现对非线性问题的有效求解。然而，隐藏层的数量和大小需要根据具体问题进行调整，过多的隐藏层可能导致模型过拟合，而过少的隐藏层可能无法捕捉到足够的特征信息。

7.7.3.2 多层感知机解决XOR问题

多层感知机能够通过其非线性映射能力解决XOR问题。在多层感知机模型中，隐藏层的主要作用是将样本从输入空间非线性映射到隐空间，进而实现特征提取和表示学习，以便更好地捕捉输入数据的复杂性和模式。

通过隐藏层的非线性变换，原本在输入空间无法线性可分的样本在隐空间中变得线性可分，从而使得模型能够准确地对XOR问题进行分类。这表明多层感知机在处理非线性问题时具有显著的优势。

7.7.4 误差反向传播算法

误差反向传播（Backpropagation）算法是一种用于训练神经网络的常见优化算法，其主要目标是通过梯度下降法调整神经网络中的参数来最小化神经网络在训练数据上的预测误差，使其能够更好地完成预测任务。

正向传播与反向传播

- **正向传播（Forward Propagation）**：在正向传播过程中，首先将输入数据传递给神经网络的输入层，并将这些输入加权求和，然后通过激活函数进行处理，最终得到输出。需注意，在正向传播过程中，每个神经元的输入是前一层神经元的输出，每一层的神经元的状态值只影响下一层神经元的状态值。
- **反向传播（Backpropagation）**：在反向传播过程中，根据整个网络的输出值与实际目标值的误差函数，计算每个参数对误差的贡献，并将这些梯度信息传播回网络，以便调整参数。从输出层开始，计算误差对每个参数的偏导数，然后逐层向前计算每层参数的梯度。计算出了参数的梯度后，使用梯度下降等优化算法更新网络中的参数，以减小误差。

通过正向传播信息和反向传播误差的重复迭代，神经网络能够不断调整参数，从而实现对数据的准确拟合和预测。然而，BP算法也存在一些问题，如梯度消失和梯度爆炸，这可能导致网络在深层结构中权重更新困难，影响模型的训练效果。

误差反向传播计算

以一个简单的神经网络为例，首先将神经网络矩阵化表示，模型输入向量 X 、权重矩阵 W 和激活函数 f （假设为Sigmoid函数）表示如下：

- 输入向量经过权重矩阵的加权和求和，可以得到隐藏层神经元的输出值 h 。
- 然后，根据误差函数计算每个参数对误差的贡献，并使用链式法则对目标函数求导，得到参数的梯度。
- 最后，使用梯度下降法更新网络中的参数，以减小误差。

通过逐层计算梯度，BP算法能够有效地将误差信息传播回网络，从而实现对参数的优化更新。然而，对于深层神经网络，梯度的计算和传播可能会受到梯度消失和梯度爆炸的影响，需要采取一些措施（如使用ReLU激活函数、批量归一化等）来缓解这些问题。

7.7.5 深度神经网络

随着数据规模和网络复杂度的增加，传统的神经网络在处理大规模数据和复杂任务时面临着挑战。为了应对这些挑战，深度神经网络应运而生。深度神经网络通过引入更多的隐藏层和参数，能更有效地学习和表示复杂的数据模式，从而提高其在大规模数据和复杂任务下的性能表现。

这种网络结构的设计使得模型能够逐层提取并组合抽象特征，逐渐建立起对输入数据更深层次、更高级别的理解。与传统神经网络相比，深度神经网络在处理图像识别、自然语言处理和语音识别等领域表现出更为优越的性能。

在通过多层的非线性变换，深度神经网络能够捕捉到数据中的复杂特征和模式，从而实现对复杂问题的准确求解。然而，深度神经网络的训练也面临着一些挑战，如计算资源需求大、训练时间长、容易过拟合等。为了克服这些挑战，研究人员提出了许多改进方法，如卷积神经网络（CNN）、循环神经网络（RNN）、长短时记忆网络（LSTM）等，这些模型在特定领域中取得了显著的成果。

7.7.6 实践案例：基于面部特征的性别分类

本节利用神经网络算法对性别进行分类，数据集由部分男女性面部特征和性别标签组成。数据预处理阶段，对离散特征进行编码，将“性别”特征转换为数值特征。模型训练时，使用Keras库构建一个全连接神经网络，并通过调整隐藏层的神经元数量和隐藏层数量来优化模型性能。

在实践案例中，我体会到了数据预处理的重要性。合理的特征编码能够使模型更好地理解 and 处理数据。同时，模型结构的选择和参数的调整也对模型性能有着重要的影响。通过增加隐藏层的数量和调整神经元数量，可以提高模型的表达能力和分类性能。然而，过多的隐藏层可能导致模型过拟合，因此需要在模型复杂度和泛化能力之间找到一个平衡点。

人工神经网络以其强大的非线性映射能力和特征提取能力，在许多实际问题中得到了广泛应用。然而，神经网络的训练也面临着一些挑战，如计算资源需求大、训练时间长、容易过拟合等。